

6/5/1

DIALOG(R)File 351:Derwent WPI

(c) 2004 Thomson Derwent. All rts. reserv.

011718460

WPI Acc No: 1998-135370/ 199813

XRFX Acc No: N98-107194

Monitoring method for system having multi-processors - by assigning predictor of capacity of system that responds quickly as evaluation object to enable allocation of processes to several processors of system whose loaded condition is monitored

Patent Assignee: MATSUSHITA ELECTRIC WORKS LTD (MATW )

Number of Countries: 001 Number of Patents: 001

Patent Family:

Patent No	Kind	Date	Applicat No	Kind	Date	Week
JP 10011407	A	19980116	JP 96164937	A	19960625	199813 B

Priority Applications (No Type Date): JP 96164937 A 19960625

Patent Details:

Patent No	Kind	Lan Pg	Main IPC	Filing Notes
JP 10011407	A		5 G06F-015/16	

Abstract (Basic): JP 10011407 A

The method involves monitoring the loaded condition of a system, that has several processors and does distributed processing, based on the priority of each process. The allocation of the processes to several processors is done by assigning the predictor of the capacity of the system that responds quickly as an evaluation object.

ADVANTAGE - Allows efficient allocation of processes to processors. Describes true behaviour of each CPU correctly. Improves overall performance and overall throughput of system.

Dwg.0/5

Title Terms: MONITOR; METHOD; SYSTEM; MULTI; PROCESSOR; ASSIGN; PREDICT; CAPACITY; SYSTEM; RESPOND; QUICK; EVALUATE; OBJECT; ENABLE; ALLOCATE; PROCESS; PROCESSOR; SYSTEM; LOAD; CONDITION; MONITOR

Index Terms/Additional Words: CENTRAL; PROCESSING; UNIT

Derwent Class: T01

International Patent Class (Main): G06F-015/16

File Segment: EPI

ref 8

(19) 日本国特許庁 (J P)

(12) 公 開 特 許 公 報 (A)

(11) 特許出願公開番号

特開平10-11407

(43) 公開日 平成10年(1998) 1月16日

(51) Int.Cl. <sup>9</sup>	識別記号	庁内整理番号	F I	技術表示箇所
G 0 6 F 15/16	3 8 0		G 0 6 F 15/16	3 8 0 Z
	3 7 0			3 7 0 N
	4 6 0			4 6 0 Z

審査請求 未請求 請求項の数 4 O L (全 5 頁)

(21) 出願番号 特願平8-164937

(22) 出願日 平成 8 年(1996) 6 月25日

(71) 出願人 000005832

松下電工株式会社

大阪府門真市大字門真1048番地

(72) 発明者 マイケル ミグドル

大阪府門真市大字門真1048番地 松下電工株式会社内

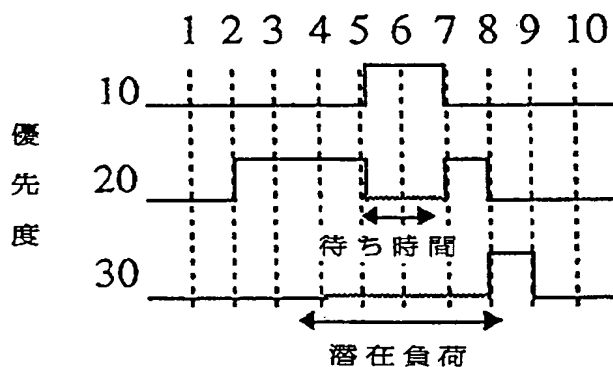
(74) 代理人 弁理士 倉田 政彦

(54) 【発明の名称】 マルチプロセッサの監視方法

(57) 【要約】

【課題】複数のプロセッサを用いて分散処理を行うシステムの負荷状態を各処理の優先度に基づいて監視し、複数のプロセッサの間で効率的に処理を分配するためのマルチプロセッサの監視方法を提供する。

【解決手段】複数のプロセッサを用いて分散処理を行うシステムの負荷状態を各処理の優先度に基づいて監視する方法において、事象に対して速やかに応答するシステムの能力の予測値、又は、プロセッサを必要とするタスク群に十分なプロセッサの時間を割り当てるシステムの能力の予測値、又は、より優先度の高いタスクによる現時点でのシステムの利用状況のうち少なくとも1つ以上に基づいて複数のプロセッサの間で効率的に処理を分配する。



## 【特許請求の範囲】

【請求項1】 複数のプロセッサを用いて分散処理を行うシステムの負荷状態を各処理の優先度に基づいて監視する方法において、事象に対して速やかに応答するシステムの能力の予測値を評価対象とすることを特徴とするマルチプロセッサの監視方法。

【請求項2】 複数のプロセッサを用いて分散処理を行うシステムの負荷状態を各処理の優先度に基づいて監視する方法において、プロセッサを必要とするタスク群に、十分なプロセッサの時間を割り当てるシステムの能力の予測値を評価対象とすることを特徴とするマルチプロセッサの監視方法。

【請求項3】 複数のプロセッサを用いて分散処理を行うシステムの負荷状態を各処理の優先度に基づいて監視する方法において、より優先度の高いタスクによる現時点でのシステムの利用状況を評価対象とすることを特徴とするマルチプロセッサの監視方法。

【請求項4】 請求項1乃至3の評価対象のうち少なくとも1つ以上に基づいて複数のプロセッサの間で効率的に処理を分配することを特徴とするマルチプロセッサの監視方法。

## 【発明の詳細な説明】

## 【0001】

【発明の属する技術分野】本発明は、複数のプロセッサを用いて分散処理を行うシステムの負荷状態を各処理の優先度に基づいて監視し、複数のプロセッサの間で効率的に処理を分配するためのマルチプロセッサの監視方法に関するものである。

## 【0002】

【従来技術】従来、優先度に基づいてオペレーティング・システムの負荷状態をモニターするための現在の方法は、一般的に生のCPUの利用率(%)又は待機しているジョブの数を用いている。ジョブの分配法を選択する方法としては、最も負荷の軽いCPUにジョブを送る方法や、各CPUに順番にジョブを送る方法、無負荷のCPUにジョブを要求する一斉同報をさせる方法を含めて、多くの方法が存在する。

【0003】これらのすべての方法の問題は、CPUの負荷の数値に依存することである。実行中のタスクの優先度が無視されているので、例えば、高優先度のタスクを実行するのに、自己の時間の80%を費やしているシステムと、低優先度のタスクを実行するのに自己の時間の80%を費やしているシステムとの間に差が生じない。その差というのは、第1の場合では、中優先度のタスクにとってはCPUの負荷が重いと見えるのに対し、第2の場合では、中優先度のタスクにとってはCPUの負荷が軽いと見えることである。このことは、与えられたジョブをどちらのCPUに実行させるべきかを選択するときに、不適切な選択を招き、その結果、全体のシステムの効率が悪くなる。

【0004】そこで、タスクの優先度に基づいて複数のプロセッサを監視する方法として、例えば、特開平7-282013号公報に開示された方法がある。

## 【0005】

【発明が解決しようとする課題】本発明の目的は、複数のプロセッサの負荷をタスクの優先度を考慮して監視する方法において、上述の従来技術には無かった3種類の評価方法を提案し、これらの方法に基づいて効率的にジョブの分配を可能とすることにある。

## 【0006】

【課題を解決するための手段】本発明のシステムでは、実行中のタスクの優先度を考慮に入れたCPUの負荷状態の新規な計測手段を開発することによって、従来例の問題点を克服している。この計測手段は、潜在負荷(事象に対して速やかに応答するシステムの能力の予測値)、タスク負荷(CPUを必要とするタスク群に、十分なCPUの時間を割り当てるシステムの能力の予測値)、又は、利用率(より優先度の高いタスクによる現時点でのシステムの利用状況)という用語で表現され得る。各CPUは、カーネルのレベルでのタスクのステータスの変化をモニターすることによって、自分自身の計測を続ける。その計測の結果は、すべてのCPUに対する、より効率的なジョブの分配を選択するために利用され得る。

【0007】以下、本発明のシステムの概略を説明する。システムの心臓部は、オペレーティング・システムのカーネルに付加された1つのソフトウェアであり、これは、タスクのステータスの変化をモニターする。次に、この情報は、3つの異なるアルゴリズムのうちの1つに伝達される。これらのアルゴリズムは、各優先レベルにおけるタスクから見た利用状況を計測し、潜在負荷を予測し、タスク負荷を予測する。各システムの個々の要求に応じて、CPUの負荷を評価するために、これらの計測手段のうちの任意の組み合わせを利用することができる。各行が優先度の数値とその優先度に関連した負荷の数値を含むデータで構成された「負荷テーブル」が出力される。

【0008】次に、CPUの情報テーブルは、CPUが密結合(例えば、同じPC基板上でISAバスを介して通信する複数のCPU)であるか、疎結合(例えば、Ethernetを介して通信する複数のCPU)であるかに応じて、2通りの方法のうちいずれかで処理される。密結合の場合、各CPUは一斉同報送信を用いた更新により自己のステータスを他のノードに知らせる。そして、各CPUは自己のタスクのうちの1つ又はそれ以上を異なるCPUに割り当てるか否かを選択することについて責任を持つ。疎結合の場合には、1つのCPUがすべての情報の管理に責任を持ち、他のCPUに対してジョブの分配方法を告げるマスターとして動作する。

## 【0009】

【発明の実施の形態】システムの動作を簡単な例により説明する。まず、第1に、1つのCPUについての負荷テーブルの計算について示す。図1は3つのタスクが実行されている簡単なケースを示す。このシステムでは、低い値の優先度は優先度の高いタスクを示している。言い換えれば、優先度10のタスクは優先度20のタスクよりも常に先にCPUを受け取る。最初に、各タスクは、休止状態であって、CPUを必要としない。時間2において、優先度20のタスクが実行可能になると、CPUを直ちに受け取る。このタスクは4ユニットの実行時間を必要とする。時間4において、優先度30のタスクが実行可能になるが、優先度のより高いタスクがCPUを使用しているので、このタスクは待機状態に入る。時間5において、優先度10のタスクが実行可能状態となる。このタスクは、優先度20のタスクよりも優先度

	1	2	3	4	5	6	7	8	9	10
10	X	X	X	X	X	O	O	X	X	X
20	X	X	O	O	O	W	W	O	X	X
30	X	X	X	X	W	W	W	W	O	X

【0012】 利用率計測のセクション  
利用率計測の計算をするために、優先レベルの各々に対するCPUの利用状況が調べられる。この例では、CPUは時間1、2及び10の間を除いて、すなわち、70%の時間は使用されていた。しかしながら、優先度30のタスクにとっては、CPUは60%の時間のみ利用不可能であった。同様に、優先度20のタスクにとっては、20%の時間だけCPUは使用され、80%の時間は利用可能であった。したがって、CPUの利用率計測モジュールの出力は、表2に示されるようなテーブルとなるであろう。優先度の値が255に対応する負荷は、生のCPUの利用率に対応していることに注意されたい。

【0013】

【表2】

優先度	CPU利用状況
10	0%
20	20%
30	60%
255	70%

【0014】 潜在負荷の予測計算のセクション  
潜在負荷を予測計算するために、或る事象が異なる時点の各々で起こったとしたときに観測されるであろう潜在負荷が計算される。例えば、優先度30のタスクが時間1又は2において実行を要求されると、このタスクは直ちにCPUを受け取るであろう。反対に、このタスクが時間5においてCPUを必要とした場合には、3ユニットの時間は待たなくてはならなかったであろう。それ故、10ユニットの時間にわたる平均的な潜在負荷は、 $(0+0+6+5+4+3+2+1+0+0)/10=$

が高いので、直ちにCPUを受け取る。優先度20のタスクは、先に実行権を得ており、まだ1ユニットの実行時間を残しているので、待機状態に入る。時間7において、優先度10のタスクが終了すると、優先度20のタスクは実行を再開し、時間8において終了する。この時点で、優先度30のタスクは、最後にCPUを受け取り、1ユニットの時間後に、その仕事を終了する。

【0010】まず、タイミングデータは、表1に示されるような単純なテーブルに変換される。このテーブルは、3つの負荷計測/性能予測のセクションの各々で使用される。テーブルの行方向は優先度、列方向は時間を表わし、は実行中、×は非実行、Wは待機中を意味する。

【0011】

【表1】

	1	2	3	4	5	6	7	8	9	10
10	X	X	X	X	X	O	O	X	X	X
20	X	X	O	O	O	W	W	O	X	X
30	X	X	X	X	W	W	W	W	O	X

2. 1ユニットの時間となる。優先度20のタスクに対する平均的な潜在負荷は、 $(0+0+0+0+0+0+2+1+0+0+0)/10=0.3$  3ユニットの時間となる。優先度10のタスクに対する平均的な潜在負荷は0である。なぜなら、CPUを使用していた、それよりも上位のタスクは存在しないから。潜在負荷の予測モジュールの出力は、表3に示されるようになるであろう。

【0015】

【表3】

優先度	予測潜在負荷
10	0
20	0.3
30	2.1
255	3.2

【0016】 待ち時間の予測計算セクション

待ち時間の予測計算のために、或る事象が異なる時点の各々で起こったとしたときに観測されるであろう待ち時間が計算される。例えば、優先度20のタスクが1、2又は7の時間でスタートすれば、待ち時間は0であろう。

他の場合には、2ユニットの待ち時間となるであろう。故に、待ち時間の予測値は、 $(0+0+2+2+2+2+0)/7=1.14$  ユニットの時間となる。時間8以降でのスタートについては、計算できないことに注意されたい。なぜなら、タスクが現在の時間とオーバーラップするであろうから。また、潜在時間は待ち時間の計算には含まれていないことに注意されたい。最後に、最も低い優先度の計算については、使用されたタスクの期間は、優先度30のタスクの時間（1ユニット）と同じであると仮定され、それ故、待ち時間は依然として0であることに注意されたい。この例の出力テーブルは、

表4のようになるであろう。

【0017】

【表4】

優先度	予測待ち時間
10	0
20	1.41
30	0
255	0

【0018】各CPUは、上に説明した情報を管理しており、周期的にそれを更新したり、古い情報を廃棄したりする。これらの数値は、2つの主な目的を有する。1つは単にCPU負荷状態のモニターとしてであり、それは、システムのパフォーマンスの異常を検出するのに有益であり得る。他の目的は、負荷均衡化システムの入力としてである。

【0019】このシステムには、2つの基本的な形態がある。1つは、CPU間の通信が低速な状況において有益であるが、そのスケジューリングの方針については正確さにおいて劣る。第2のものは、高速なプロセッサ間通信が可能となるときに有益であり、プロセッサ間で交換される情報量が増加することにより、より正確なものとなる。以下、それぞれの場合の実施例について説明する。

【0020】

【実施例】

(A) 低速法の実施例

この方法では、その全体的な負荷が所定のスレシールドよりも少ないプロセッサ群のみが、それらの負荷テーブル(表2に示されるような)を他のすべてのノードに一斉同報送信する。そのとき、負荷の重いノードは、自己

が知る無負荷のノードについての情報に従って、自己のタスク群の1つを引き継ぐことをそのノードに要求する。

【0021】一例として、図2に示される簡単な例を考える。この例は、3つのノードを有しており、負荷評価システムは、CPU利用率のみを用いる。この例では、ノード1は、その上で実行される3つのタスクを有している。高優先度のタスクはCPUの時間の20%を取っており、中優先度のタスクはCPUの時間の10%を取っており、低優先度のタスクはCPUの時間の60%を取っている。ノード2及び3の負荷は軽く、それ故、これら両者は、現時点の自己の利用率テーブルを一斉同報送信する。ノード1は中優先度のジョブを他のノードに送りたい。なぜなら、その高優先度のジョブが中優先度のジョブについて緩慢な応答を引き起こしているからである。従来の負荷均衡化システムでは、全体の利用率の値が小さい方のノード2にジョブを送る傾向があった。このことは、実際には、中優先度のノードのパフォーマンスの低下を招いている。なぜなら、ノード2の高優先度のプロセスはCPUの時間の30%を取っているから

である。しかし、中優先度のジョブは、ノード3上では、もっと良い応答を得られるであろうということをノード1は理解している。そして、それ故に、図3に示されるように、ジョブをノード3に送る。現実の利用分布は多く変更されていないけれども、中優先度のタスクのパフォーマンスは、ノード3にシフトすることによって、かなりの程度、改善されるに違いないことに注意されたい。

【0022】(B) 高速法の実施例

この方法では、各プロセッサは、自己のタイミング情報を(負荷測定情報の代わりに)マスタープロセッサに送る。そのタイミング情報は、(表1に示されるような)タイミングテーブルを含み、また、将来のプロセッサのタスク群についてのタイミング要求について既知の任意の情報をも含んでいる。このマスタープロセッサは、送られてきたタイミングプロフィールを検討して、自己が受信した情報に基づいて、すべてのプロセッサのために最適のジョブの配置を決定する。タイミング情報を送信するのに必要な帯域幅のために、高速のプロセッサ間通信が必要である。この結果、より効率的な配置を選択することができ、全体的なパフォーマンスを改善できる。

【0023】この方法は、図4に示される。ノード1はマスターのノードであり、自分自身のタイミングデータを処理すると共に、ノード2及び3から送られてきたタイミングデータも処理する。マスターが一度データを処理すると、図5に示されるように、マスターは選択された配置を一斉同報送信し、そのとき、各ノードは、マスターのノードによって選択されたように、ジョブを引き継ぐ。

30 【0024】

【発明の効果】このシステムでは、負荷計算において実行中のタスクの優先度を用いることによって、生のCPUの利用と優先度に基づくCPUの利用との間に明瞭な区別をつけることができ、より正確に各CPUの真の振る舞いを記述できる。その結果、複数のタスクは、それらの優先度に基づいて、より正確に分配され、システムのパフォーマンスと全体的なスループットを改善できる。

【図面の簡単な説明】

40 【図1】本発明の監視方法によるマルチタスクの監視状況を示す説明図である。

【図2】本発明の監視方法を低速なプロセッサ通信によるマルチプロセッサシステムに適用した場合のジョブの再分配前の状態を示す説明図である。

【図3】本発明の監視方法による低速なプロセッサ通信によるマルチプロセッサシステムに適用した場合のジョブの再分配後の状態を示す説明図である。

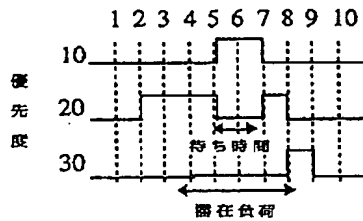
50 【図4】本発明の監視方法を高速なプロセッサ通信によるマルチプロセッサシステムに適用した場合のマスタープロセッサへの情報の流れを示す説明図である。

【図5】本発明の監視方法を高速なプロセッサ通信によるマルチプロセッサシステムに適用した場合のマスタープロセッサからの情報の流れを示す説明図である。

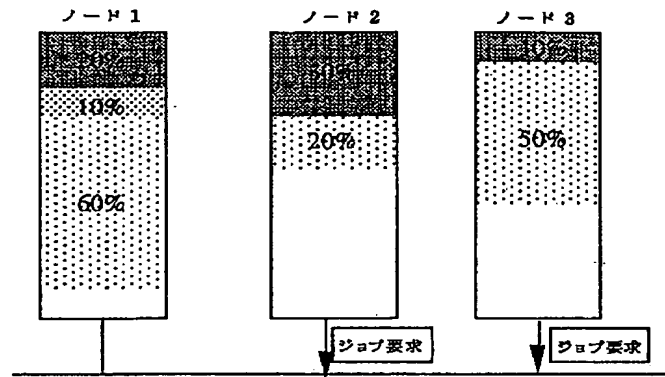
【符号の説明】

1～3 ノード

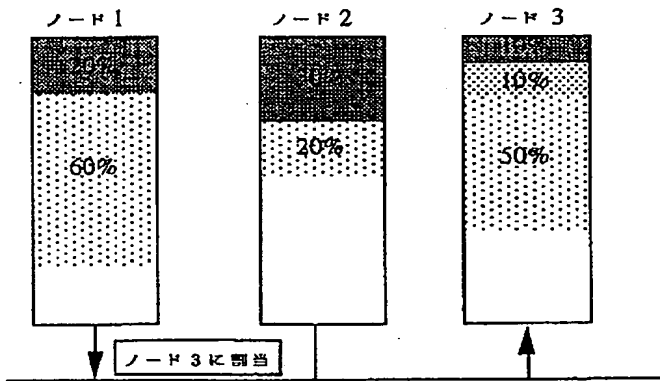
【図1】



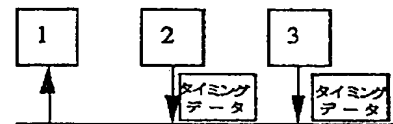
【図2】



【図3】



【図4】



【図5】

